# NUS-IDS at PragTag-2023: Improving Pragmatic Tagging of Peer Reviews through Unlabeled Data

**Sujatha Das Gollapalli**
Institute of Data Science
National University of Singapore
Singapore
idssdg@nus.edu.sg

**Yixin Huang***
Télécom Paris
Institut Polytechnique
de Paris, France
yixin.huang@ip-paris.fr

**See-Kiong Ng**
Institute of Data Science
National University of Singapore
Singapore
seekiong@nus.edu.sg

## Abstract

We describe our models for the *Pragmatic Tagging of Peer Reviews* Shared Task at the 10th Workshop on Argument Mining at EMNLP-2023. We trained multiple sentence classification models for the above competition task by employing various state-of-the-art transformer models that can be fine-tuned either in the traditional way or through instruction-based fine-tuning. Multiple model predictions on unlabeled data are combined to tentatively label unlabeled instances and augment the dataset to further improve performance on the prediction task. In particular, on the F1000RD corpus, we perform on-par with models trained on 100% of the training data while using only 10% of the data. Overall, on the competition datasets, we rank among the top-2 performers for the different *data conditions*.

## 1 Introduction

**Peer Review** is employed across various subject domains to assess the quality of research documents such as grant proposals, journal manuscripts, and conference proceedings. Peer reviews are performed by independent researchers with expertise on the relevant topic for purposes such as awarding grants or publishing latest research for the advancement of Science. Review text reports, the result of these peer assessments, are brief summaries describing the document's main contributions, its strengths and weaknesses, along with other revision related comments and constructive feedback (Griessenauer and Roach, 2019).

Though standards and practices may vary across different subject domains and even across venues within the same domain, the main objective of the peer review process is to ensure the advancement of quality research (Glonti et al., 2019). To this end, alleviating the reviewing burden and supporting the diverse nature of reviewer expertise becomes vital (Huisman and Smits, 2017) and motivates the on-going research on developing tools to assist and improve the peer reviewing process (Walker and Rocha da Silva, 2015; Checco et al., 2021; Yuan et al., 2022; Schulz et al., 2022). In particular, a significant direction towards developing AI-assisted peer reviewing models involves the compilation of relevant datasets to support the meta-analyses of reviews (Kang et al., 2018; Ghosal et al., 2022; Dycke et al., 2023a).

From the perspective of language and NLP research, review reports provide a rich ground for investigation for various argument mining problems (Hua et al., 2019) including classification tasks such as paper acceptance prediction and sentence labeling (Bao et al., 2021; Kuznetsov et al., 2022). The PragTag Shared Task[1] at the 10th Workshop on Argument Mining at EMNLP-2023 comprises one such sentence labeling task in which every sentence from a review report is assigned a label from one of the pragmatic categories: { Recap, Strength, Weakness, Todo, Other, Structure}. Due to space constraints, we refer our readers to Kuznetsov, et al. (2022) and Dycke, et al. (2023b) for the precise definitions of the pragmatic categories and the F1000RD Corpus which forms the basis for the datasets used in the PragTag-2023 competition.

### 1.1 Task Description and Evaluation

In PragTag-2023, the pragmatic tagging task is presented in a cross-domain, low-resource setting using data from the F1000RD Corpus. The F1000RD is a multi-domain collection of free-text peer reviews annotated with pragmatic labels at the sentence level. Each peer review is associated with a domain (related to Medicine, Computer Science, or Scientific Policy Research). Additionally, recently released unlabeled review corpora from Dycke, et

---

[1]https://codalab.lisn.upsaclay.fr/competitions/13334

al. (2023a) were made available as auxiliary data sources. The following three **data conditions** were proposed for the competition:

1. No-data: where no labeled instances are available for the task–zero-shot setting (Radford et al., 2019).

2. Low-data: where about 20% of the labeled data for the task can be used for training models–few-shot setting (Brown et al., 2020).

3. Full-data: which is the standard machine-learning setting where the entire training split of the labeled data can be used to train models.

For measuring model performance on this sentence classification task, the average performance across domains is used in each of the above conditions where the performance in a domain is simply measured by the macro-F1 computed across all review sentences of that domain. For the final evaluation, the test data comes from a "secret" domain, different from those covered in the training data, thus measuring cross-domain model performance.

---

Consider the definitions of labels:
*Recap: summarizes the manuscript, For e.g. "The paper proposes a new method for...";*
... Question: Which of the above labels most applies to the following sentence? Sentence: []

---

Table 1: Prompt for LLM Models

## 2 Proposed Methods

In this section, we briefly describe the various models we employed for the Pragmatic Tagging task under the three data conditions.

**No-data setting**: We studied two approaches for predicting pragmatic tags under the no-data condition. In the first "*Semantic Search*" approach, we simply use a list of "questions" to find sentences in the review texts that best answer the question. This list was curated based on the typical questions employed during the peer review process of NLP conferences and augmented to cover labels such as "Recap".[2] Example questions include "How original are the results described in the paper?" and "What is the main finding of this paper?". We used the state-of-the-art Sentence Transformer models

---

[2]Complete list shared as part of the code distribution

trained for Semantic Search for this method (Wang et al., 2020; Nassiri and Akhloufi, 2023).[3]

Recent breakthrough research has shown that large language models (LLMs) can be trained "to act in accordance with the user's intentions" and as a consequence be "prompted" to perform a range of NLP tasks (Radford et al., 2019; Brown et al., 2020; Christiano et al., 2017). For our second approach, in keeping with this recent direction, we designed a multiple-choice question prompt along with the task description provided in the competition for use in Instruction Fine-tuned Language Models (Ouyang et al., 2022; Chung et al., 2022). Our prompt is listed in Table 1 and we refer to the use of this approach as "*MC-Prompt*" in Section 3.[4]

**Low-data/Full-data setting**: In current practice, fine-tuning large pre-trained language models (PLMs) for a new task has become the standard approach for training models (Howard and Ruder, 2018). We therefore adopt the state-of-the-art transformer-based models and directly train supervised models on the available labeled data for the low/full data conditions.

With the objective to utilize the unlabeled data provided in the competition as means to overcome the scarcity of labeled data in the low-data settings, we employed traditional semi-supervised approaches–self-training and voting, to combine predictions from multiple learners[5] and obtain tentative labels for the unlabeled data (Li et al., 2019; Sosea and Caragea, 2022). The "tentatively labeled" unlabeled data is incorporated via two methods in our models. In the pretraining approach (*PT*), we simply pretrain our classifier on the tentatively-labeled unlabeled data before fine-tuning on the labeled data whereas in the *Combined* approach, the augmented dataset is used to train a model.

## 3 Experiments

**Datasets**: We used the datasets from previous works (Kuznetsov et al., 2022; Dycke et al., 2023a) for showcasing our proposed methods on this task.

---

[3]https://huggingface.co/sentence-transformers/multi-qa-mpnet-base-dot-v1

[4]We experimented with slight variations and paraphrases of the label descriptions, prompts with and without examples, as well as a yes/no prompt that uses a yes/no question with each label. Our best prompt based on validation performance is listed in Table 1.

[5]In addition to the provided RoBERTa-based competition baseline, we also fine-tuned models based on T5 and FlanT5 models from Google. These details are provided in Section 3.

| Setting | #Labeled Sentences | Model | Accuracy | Macro-F1 |
|---|---|---|---|---|
| No data | 0 | QA-MPNet (*Semantic Search*) | 0.31 | 0.32 |
| | 0 | FlanT5-XL (*MC-Prompt*) | **0.53** | **0.42** |
| Low data (10%) | 338 | RoBERTa | 0.75 | 0.71 |
| | 338 | FlanT5-large | 0.70 | 0.68 |
| | 338 | T5-large | 0.78 | 0.74 |
| | 14673/338 | T5-large (*PT*) | **0.81** | **0.80** |
| | 15011 | T5-large (*Combined*) | 0.77 | 0.76 |
| Full-data | 2691 | Roberta | 0.83 | 0.82 |
| | 2691 | FlanT5-large | 0.82 | 0.80 |
| | 2691 | T5-large | 0.84 | 0.82 |
| | 15844/2691 | T5-large (*PT*) | **0.86** | **0.85** |
| | 18535 | T5-large (*Combined*) | 0.85 | 0.83 |

Table 2: Performance of various models is shown on the test split of the F1000RD corpus. The best performance in each setting is highlighted in bold. For the "X/Y" values shown in the #Labeled Sentences column of PT rows, X is the number of tentatively-labeled unlabeled instances and Y, the number of labeled instances from the training data.

In particular, we used the F1000RD Corpus[6] for presenting our observations in this section. For the competition, in accordance with the competition rules, we only used the provided main and auxiliary datasets (Dycke et al., 2023b).

**Implementation Details**: We fine-tuned the Text-to-Text Transfer Transformer (T5) model for our classification task. T5 incorporates various tasks such as translation, question answering, and classification uniformly as text-to-text learning tasks, thereby harnessing the power of transfer learning across multiple tasks, and has been shown to obtain state-of-the-art performance across a range of tasks (Raffel et al., 2020). The T5 models were extended to incorporate instruction-based fine-tuning into the FlanT5-family of models (Chung et al., 2022). For T5 and FlanT5 experiments, we used latest implementations available from Hugging-Face (Wolf et al., 2019). In total, for the low/full data conditions, three classifiers were trained using T5, FlanT5, and the RoBERTa baseline provided in the competition.

All experiments were performed on a single GPU of an Nvidia Tesla cluster machine with 32GB RAM. On this machine, based on the size of the datasets and the specific models, training time ranges between 0.5-24 hours. On our available infrastructure, the biggest models we were able to train were the "large" variants (T5-large and FlanT5-large) from the T5 and FlanT5

model families. The performance on the development/validation split of the dataset was used to set the number of epochs for the final models.[7]

## 3.1 Results and Observations

We illustrate the performance of our models under the three data conditions on the F1000RD dataset.[6] For the *low-data* condition shown in Table 2, we used a randomly-selected 10% subset of the training data. In this table, we see that, not surprisingly, the accuracy and macro-F1 scores of models in the full-data condition are significantly higher than those in the low-data condition. However, in absolute terms, even with 10% of the labeled data the performance is reasonably high on this dataset. Moreover, using appropriate prompts in the FlanT5-XL model, we are able to obtain almost half of the Macro-F1 score obtained with full-data models even in the no-data condition.

Based on the competitive validation performance afforded by the T5-large models in both *low-data* and *full-data* conditions, we selected this model for exploring the improvements with unlabeled data. For these two data conditions, we used the three models (RoBERTa, FlanT5-large, T5-large) to obtain predictions for the auxiliary (unlabeled) data made available in the competition. We incorporate those examples for which there is agreement between RoBERTa and FlanT5-large model predictions but no agreement with T5-large model

| Setting | Model | F1-case | F1-diso | F1-rpkg | F1-iscb | F1-scip | F1-mean |
|---|---|---|---|---|---|---|---|
| No-data | QA-MPNet | 0.352 | 0.310 | 0.354 | 0.326 | 0.291 | 0.326 |
| | FlanT5-large | 0.420 | 0.396 | 0.413 | 0.424 | 0.357 | 0.402* |
| **Rank-1** | Unknown | 0.502 | 0.518 | 0.492 | 0.551 | 0.516 | 0.516 |
| Low-data | T5-large | 0.764 | 0.792 | 0.789 | 0.796 | 0.827 | 0.794 |
| (**Rank-1**=Us) | FlanT5-large | 0.804 | 0.835 | 0.803 | 0.803 | 0.820 | 0.813* |
| Full-data | T5-large | 0.813 | 0.853 | 0.829 | 0.806 | 0.861 | 0.832 |
| | T5-large (*PT*) | 0.843 | 0.834 | 0.827 | 0.821 | 0.854 | 0.836 |
| | T5-large (*Combined*) | 0.838 | 0.854 | 0.848 | 0.833 | 0.878 | 0.850* |
| **Rank-1** | Unknown | 0.829 | 0.842 | 0.854 | 0.836 | 0.889 | 0.850 |

Table 3: Phase-1 Results from the competition. We indicate the performance of the best system in the **Rank-1** row and highlight our best F1-mean score with a *

| Setting | Model | F1-secret | F1-mean |
|---|---|---|---|
| No | FlanT5-large | 0.425 | 0.406 |
| Low | FlanT5-large | 0.759 | 0.804 |
| Full | T5-large (*Combined*) | 0.741 | 0.832 |
| Rank-1 | Unknown | 0.801 | 0.841 |

Table 4: Phase-2 Results. The Rank-1 row shows the performance of the best model from the competition.

| Class Label | Default | PT | Combined |
|---|---|---|---|
| Other | 0.63 | **0.70** | 0.62 |
| Recap | 0.74 | **0.80** | *0.77 |
| Strength | 0.83 | **0.87** | *0.85 |
| Structure | 0.95 | 0.92 | **0.95** |
| Todo | 0.94 | **0.95** | 0.94 |
| Weakness | 0.84 | **0.85** | ***0.85** |
| Macro Average | 0.82 | **0.85** | 0.83 |

Table 5: Test F1 performance for each class label is shown for the three T5-large models from Table 2. The best performances are **bolded**. We also highlight the cases where the *Combined* setting outperforms the default setting with a *.

predictions as the subset of "weakly-labeled" data for training new T5 models in *PT* and *Combined* settings described in Section 2.

That is, during data augmentation, we add the "hard" cases for which the T5-large model predictions do not match the labels predicted by both RoBERTa and Flan-T5. This step cuts down the amount of unlabeled data added back to the dataset by excluding "uninformative" samples for which the original T5 model predictions already conform to the other models. In our early experiments, we observed that adding all examples for which we have majority labels significantly increases the training time with no significant improvements in the validation performance.

As can be seen in Table 2, both *PT* and *Combined* settings result in improved test performance for *low-data* as well as the *full-data* conditions. In particular, the improvement is significantly higher in the macro F1 score in the *low-data* condition. Indeed, with pretraining (*PT*), the test performance in low-data conditions is comparable to those of models trained on full data.

In Table 5, the per-class F1 scores on the test split for the three models: T5-large, T5-large (PT), T5-large (Combined) from Table 2 are shown. The improved F1 scores across classes in both *PT* and *Combined* settings are indicative of a significant reduction in the number of erroneous predictions

over the baseline setting. As such, F1 improvements are seen for five out of the six classes in the *PT* setting, and three out of the six classes in the *Combined* setting.

## 3.2 Competition Performance and Ranking

The results with our models in the competition are showcased for the two phases in Tables 3 and 4. Within the competition timeframe and limits on number of submissions, we were unable to test all our models on the final dataset. We highlight our best-performing models among those we submitted and also the overall best submission in the competition (Rank 1) for each condition. During the competition, for the *PT* and *Combined* runs, we used all unlabeled examples with majority labels (different from the settings used in Table 2).

Overall, we ranked among the top-2 performing of the four-six submitted systems for the various data conditions. Compared to the performances highlighted in Tables 2 and 3, our models underperform on the data from the secret domain (Table 4) indicating that they may not be generalizing well for new/unseen domains.

## 4 Related Work

Sentence classification tasks are well-studied in NLP research with deep learning models comprising the state-of-the-art (Cohan et al., 2019). Some recent sentence-level classification tasks include identification of complex linguistic phenomena in texts such as emotions, empathy, humor, sarcasm, and dialog acts (Song et al., 2022; He et al., 2021; Wang et al., 2022; Bunescu and Uduehi, 2022).

Recently, efforts are underway for collecting relevant datasets for designing assistive automation aids for peer review (Yuan et al., 2022; Checco et al., 2021; Kang et al., 2018; Ghosal et al., 2022; Dycke et al., 2023a). In this context, Kuznetsov, et al. (2022) introduced pragmatic tagging for labeling sentences of peer reviews using a schema that applies across different research fields and communities. We borrow from the latest NLP advances such as prompt-based models and combine them with unlabeled data on precisely this task.

## 5 Conclusions and Future Work

We presented our approaches for the pragmatic tag prediction task for peer reviews as part of the Prag-Tag Shared Task @ ArgMining Workshop 2023. In particular, we studied prompt-based fine-tuning as a viable alternative to traditional learning methods for this task and showcased how unlabeled data may be utilized via multiple learners to improve performance in the low-data settings. In future, we would like to address the generalizability of our proposed models across various subject domains as well as extend our approaches to related tasks such as paper acceptance prediction (Bao et al., 2021; Yuan et al., 2022).

## Acknowledgments

## References

Peng Bao, Weihui Hong, and Xuanya Li. 2021. Predicting paper acceptance via interpretable decision sets. In *Companion Proceedings of the Web Conference 2021*, WWW '21, page 461–467, New York, NY, USA. Association for Computing Machinery.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Razvan C. Bunescu and Oseremen O. Uduehi. 2022. Distribution-based measures of surprise for creative language: Experiments with humor and metaphor. In *Proceedings of the 3rd Workshop on Figurative Language Processing (FLP)*, pages 68–78, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Alessandro Checco, Lorenzo Bracciale, Pierpaolo Loreti, Stephen Pinfield, and Giuseppe Bianchi. 2021. Ai-assisted peer review. *Humanities and Social Sciences Communications*, 8.

Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Arman Cohan, Iz Beltagy, Daniel King, Bhavana Dalvi, and Dan Weld. 2019. Pretrained language models for sequential sentence classification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3693–3699, Hong Kong, China. Association for Computational Linguistics.

Nils Dycke, Ilia Kuznetsov, and Iryna Gurevych. 2023a. NLPeer: A unified resource for the computational study of peer review. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5049–5073, Toronto, Canada. Association for Computational Linguistics.

Nils Dycke, Ilia Kuznetsov, and Iryna Gurevych. 2023b. Overview of PragTag-2023: Low-resource multi-domain pragmatic tagging of peer reviews. In *Proceedings of the 10th Workshop on Argument Mining*, Singapore. Association for Computational Linguistics.

Tirthankar Ghosal, Kamal Kaushik Varanasi, and Valia Kordoni. 2022. Hedgepeer: A dataset for uncertainty detection in peer reviews. In *Proceedings of the 22nd ACM/IEEE Joint Conference on Digital Libraries*, JCDL '22, New York, NY, USA. Association for Computing Machinery.

Ketevan Glonti, Daniel Cauchi, Erik Cobo, Isabelle Boutron, David Moher, and Darko Hren. 2019. A scoping review on the roles and tasks of peer reviewers in the manuscript review process in biomedical journals. *BMC medicine*, 17:1–14.

Christoph J Griessenauer and Michelle K Roach. 2019. Scientific peer review. *A Guide to the Scientific Career: Virtues, Communication, Research and Academic Writing*, pages 403–406.

Zihao He, Leili Tavabi, Kristina Lerman, and Mohammad Soleymani. 2021. Speaker turn modeling for dialogue act classification. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2150–2157, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In *ACL*, pages 328–339.

Xinyu Hua, Mitko Nikolov, Nikhil Badugu, and Lu Wang. 2019. Argument mining for understanding peer reviews. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2131–2137, Minneapolis, Minnesota. Association for Computational Linguistics.

Janine Huisman and Jeroen Smits. 2017. Duration and quality of the peer review process: the author's perspective. *Scientometrics*, 113(1):633–650.

Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. 2018. A dataset of peer reviews (PeerRead): Collection, insights and NLP applications. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1647–1661, New Orleans, Louisiana. Association for Computational Linguistics.

Ilia Kuznetsov, Jan Buchmann, Max Eichler, and Iryna Gurevych. 2022. Revise and resubmit: An intertextual model of text-based collaboration in peer review. *Computational Linguistics*, 48(4):949–986.

Xinzhe Li, Qianru Sun, Yaoyao Liu, Shibao Zheng, Qin Zhou, Tat-Seng Chua, and Bernt Schiele. 2019. *Learning to Self-Train for Semi-Supervised Few-Shot Classification*.

Khalid Nassiri and Moulay Akhloufi. 2023. Transformer models used for text-based question answering systems. *Applied Intelligence*, 53(9):10602–10635.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21(140):1–67.

Robert Schulz, Adrian Barnett, René Bernard, Nicholas J.L. Brown, Jennifer A. Byrne, Peter Eckmann, Małgorzata A. Gazda, Halil Kilicoglu, Eric M. Prager, Maia Salholz-Hillel, Gerben ter Riet, Timothy Vines, Colby J. Vorland, Han Zhuang, Anita Bandrowski, and Tracey L. Weissgerber. 2022. Is the future of peer review automated? *BMC Research Notes*, 15(1).

Xiaohui Song, Longtao Huang, Hui Xue, and Songlin Hu. 2022. Supervised prototypical contrastive learning for emotion recognition in conversation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5197–5206, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Tiberiu Sosea and Cornelia Caragea. 2022. Leveraging training dynamics and self-training for text classification. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4750–4762, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Richard Walker and Pascal Rocha da Silva. 2015. Emerging trends in peer review—a survey. *Frontiers in Neuroscience*, 9.

Jiquan Wang, Lin Sun, Yi Liu, Meizhi Shao, and Zengwei Zheng. 2022. Multimodal sarcasm target identification in tweets. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8164–8175, Dublin, Ireland. Association for Computational Linguistics.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA. Curran Associates Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771.

Weizhe Yuan, Pengfei Liu, and Graham Neubig. 2022. Can we automate scientific reviewing? *J. Artif. Int. Res.*, 75.